

ЛАБОРАТОРНА РОБОТА № 4

ДОСЛІДЖЕННЯ МЕТОДІВ АНСАМБЛЕВОГО НАВЧАННЯ ТА СТВОРЕННЯ РЕКОМЕНДАЦІЙНИХ СИСТЕМ

Мета роботи: використовуючи спеціалізовані бібліотеки та мову програмування Python дослідити методи ансамблів у машинному навчанні та створити рекомендаційні системи.

Завдання 4.1 Створення класифікаторів на основі випадкових та гранично випадкових лісів.

```
1 import argparse
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from sklearn.metrics import classification_report
5 from sklearn.model_selection import train_test_split
6 from sklearn.ensemble import RandomForestClassifier, ExtraTreesClassifier
7 from sklearn.metrics import classification_report
8 from utilities import visualize_classifier
9 # Парсер аргументів
10 def build_arg_parser():
11     parser = argparse.ArgumentParser(description='Classify data using Ensemble Learning techniques')
12     parser.add_argument('--classifier-type', dest='classifier_type',
13                         required=True, choices=['rf', 'ertf'],
14                         help='Type of classifier to use; can be either 'rf' or 'ertf')
15     return parser
16 if __name__ == '__main__':
17     # Вилучення вхідних аргументів
18     args = build_arg_parser().parse_args()
19     classifier_type = args.classifier_type
20     # Завантаження вхідних даних
21     input_file = 'data_random_forests.txt'
22     data = np.loadtxt(input_file, delimiter=',')
23     X, y = data[:, :-1], data[:, -1]
24     # Розбиття вхідних даних на три класи
25     class_0 = np.array(X[y == 0])
26     class_1 = np.array(X[y == 1])
27     class_2 = np.array(X[y == 2])
28     # Візуалізація вхідних даних
29     plt.figure()
30     plt.scatter(class_0[:, 0], class_0[:, 1], s=75, facecolors='white',
31                 edgecolor='black', linewidth=1, marker='s')
32     plt.scatter(class_1[:, 0], class_1[:, 1], s=75, facecolors='white',
33                 edgecolor='black', linewidth=1, marker='o')
34     plt.scatter(class_2[:, 0], class_2[:, 1], s=75, facecolors='white',
35                 edgecolor='black', linewidth=1, marker='^')
36     plt.title('Input data')
37     # Розбивка даних на навчальний та тестовий набори
38     X_train, X_test, y_train, y_test = train_test_split(
39         X, y, test_size=0.25, random_state=5)
40     # Класифікатор на основі ансамблевого навчання
41     params = {'n_estimators': 100, 'max_depth': 4, 'random_state': 0}
42     if classifier_type == 'rf':
43         classifier = RandomForestClassifier(**params)
44     else:
45         classifier = ExtraTreesClassifier(**params)
46     classifier.fit(X_train, y_train)
47     visualize_classifier(classifier, X_train, y_train, 'training dataset')
48     y_test_pred = classifier.predict(X_test)
49     visualize_classifier(classifier, X_test, y_test, 'test dataset')
50     # Перевірка роботи класифікатора
51     class_names = ['Class-0', 'Class-1', 'Class-2']
52     print("\n" + "#" * 40)
53     print("\nClassifier performance on training dataset\n")
54     print(classification_report(y_train, classifier.predict(X_train), target_names=class_names))
55     print("#" * 40 + "\n")
56     print("\nClassifier performance on test dataset\n")
57     print(classification_report(y_test, y_test_pred, target_names=class_names))
58     print("#" * 40 + "\n")
59     # Виведення параметрів класифікатора
```

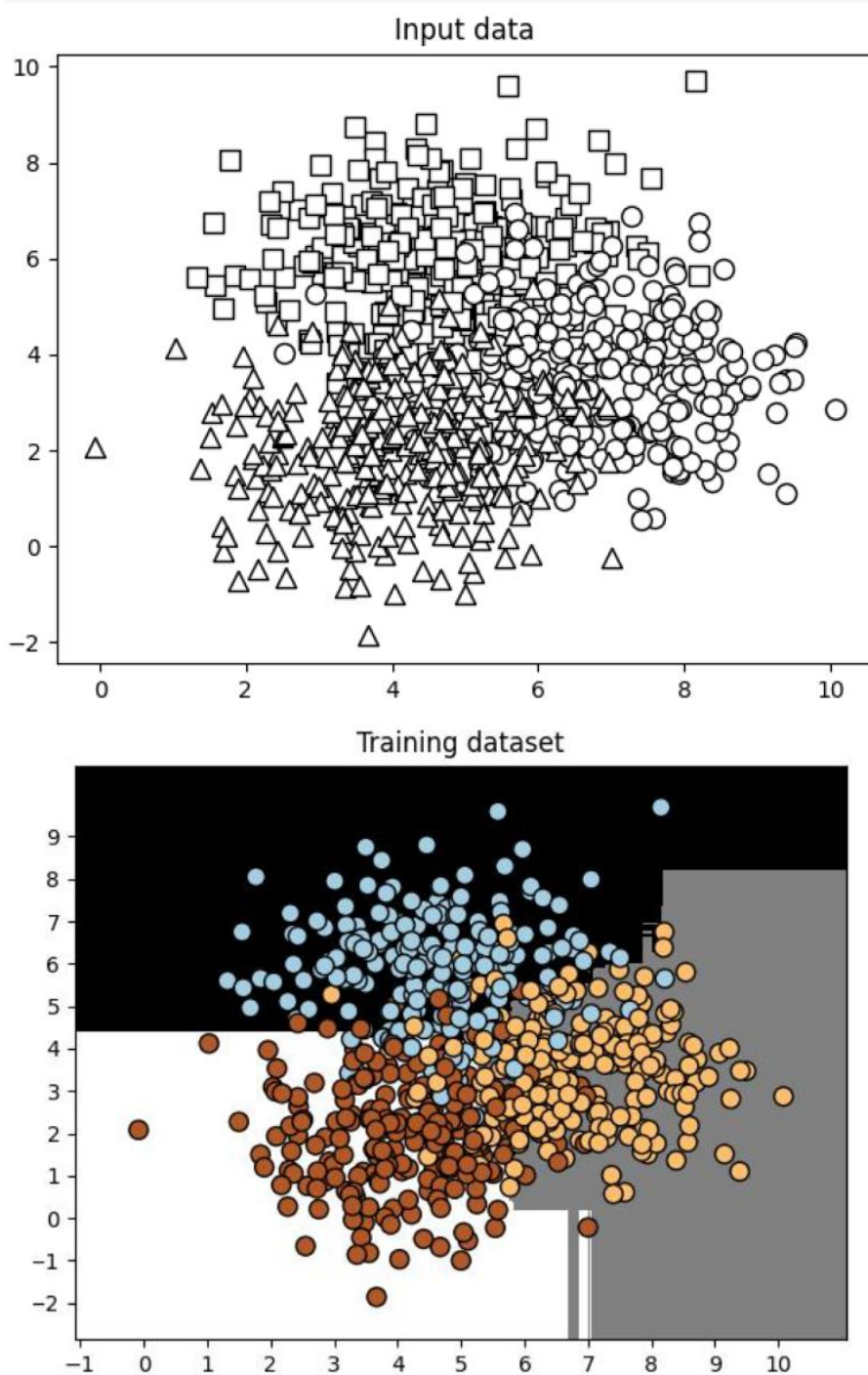
					ДУ «Житомирська політехніка».20.121.12.			
Змн.	Арк.	№ докум.	Підпис	Дата				
Розроб.		Надворний М.Ю.				Літ.	Арк.	Аркушів
Перевір.		Філіпов В.О.					1	24
Керівник						ФІКТ Гр. ІПЗк-20-1		
Н. контр.								
Зав. каф.								

```

#Обчислення параметрів довірливості
test_datapoints = np.array([[5, 5], [3, 6], [6, 4], [7, 2], [4, 4], [5, 2]])
print("\nConfidence measure:")
for datapoint in test_datapoints:
    probabilities = classifier.predict_proba([datapoint])[0]
    predicted_class = 'Class-' + str(np.argmax(probabilities))
    print('\nDatapoint:', datapoint)
    print('Predicted class:', predicted_class)
#Візуалізація точок даних
visualize_classifier(classifier, test_datapoints, [0] * len(test_datapoints), 'Test datapoints')
plt.show()

```

Рис 1. Код програми файлу LR_4_task_1

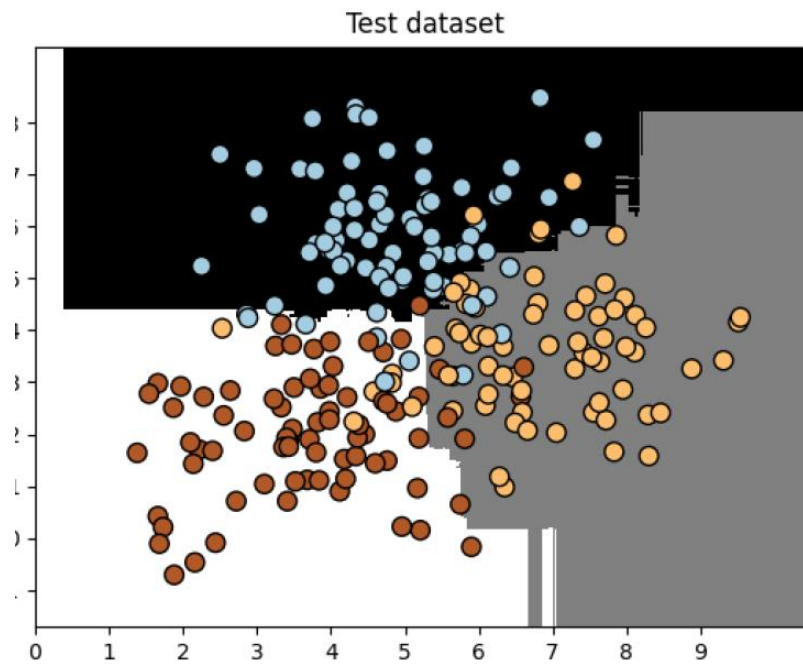


		Надворний М.Ю		
		Філіпов В.О.		
Змн.	Арк.	№ докум.	Підпис	Дата

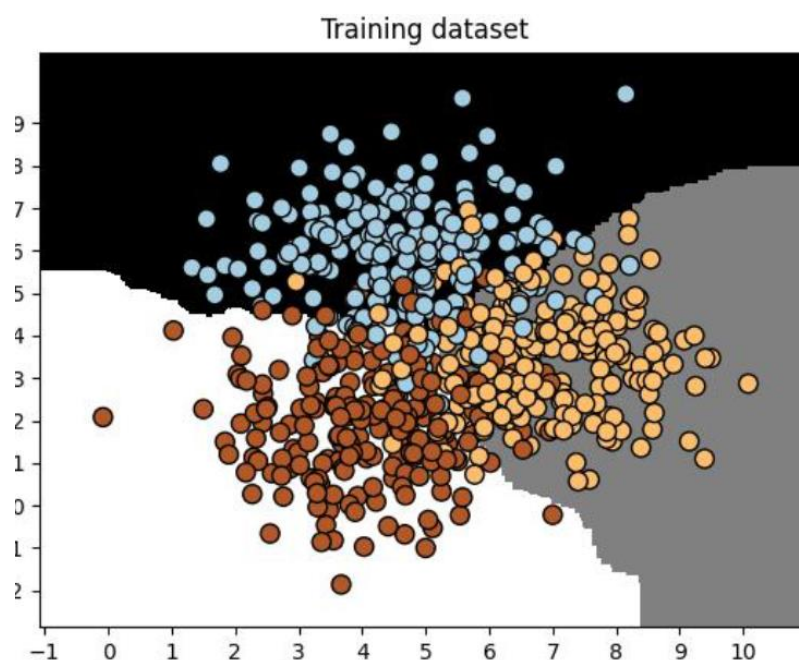
ДУ «Житомирська політехніка».20.121.12 – Лр4

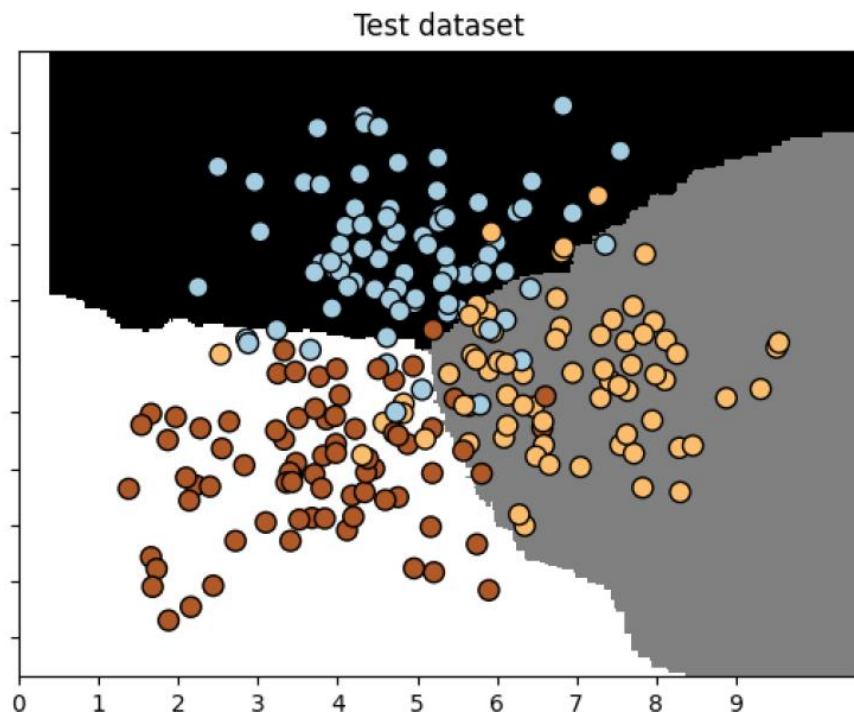
Арк.

2



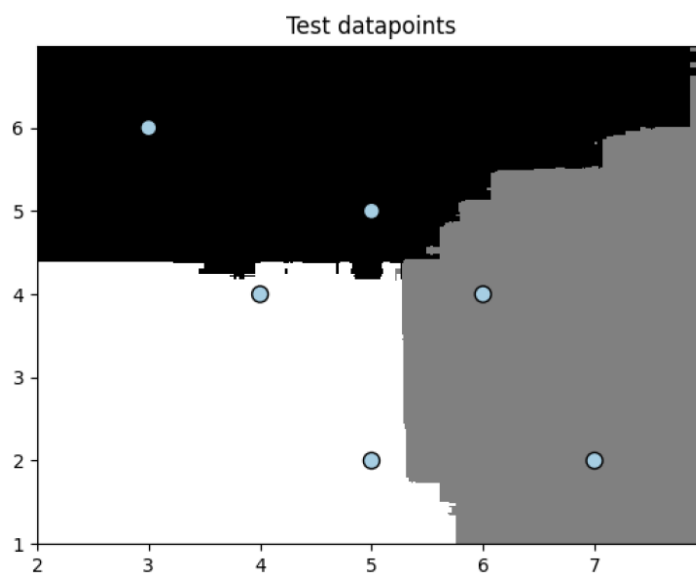
Class-0	0.92	0.85	0.88	79
Class-1	0.86	0.84	0.85	70
Class-2	0.84	0.92	0.88	76
accuracy			0.87	225
macro avg	0.87	0.87	0.87	225
weighted avg	0.87	0.87	0.87	225
#####				





Classifier performance on test dataset

	precision	recall	f1-score	support
Class-0	0.92	0.85	0.88	79
Class-1	0.84	0.84	0.84	70
Class-2	0.85	0.92	0.89	76
accuracy			0.87	225
macro avg	0.87	0.87	0.87	225
weighted avg	0.87	0.87	0.87	225



```

Confidence measure:

Datapoint: [5 5]
Predicted class: Class-0

Datapoint: [3 6]
Predicted class: Class-0

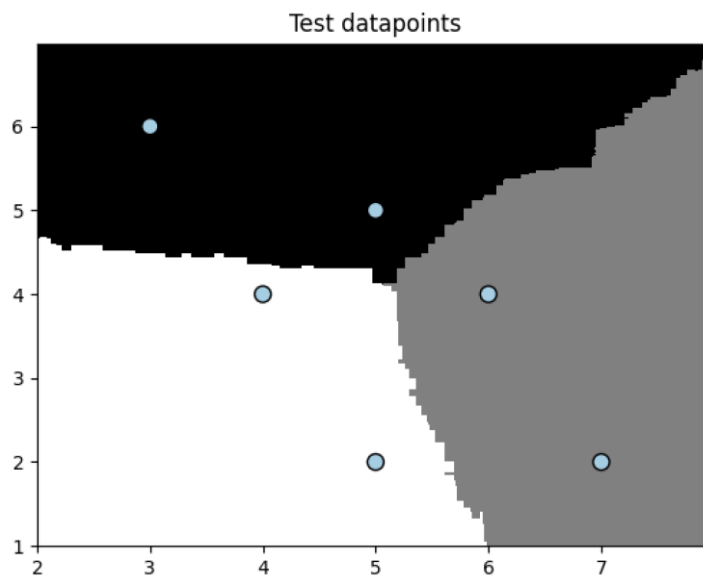
Datapoint: [6 4]
Predicted class: Class-1

Datapoint: [7 2]
Predicted class: Class-1

Datapoint: [4 4]
Predicted class: Class-2

Datapoint: [5 2]
Predicted class: Class-2

```



		Надворний М.Ю			ДУ «Житомирська політехніка».20.121.12 – Лр4	Арк.
		Філіпов В.О.				
Змн.	Арк.	№ докум.	Підпис	Дата		5

```

Confidence measure:

Datapoint: [5 5]
Predicted class: Class-0

Datapoint: [3 6]
Predicted class: Class-0

Datapoint: [6 4]
Predicted class: Class-1

Datapoint: [7 2]
Predicted class: Class-1

Datapoint: [4 4]
Predicted class: Class-2

Datapoint: [5 2]
Predicted class: Class-2

```

Рис 2. Результат виконання коду файлу LR_4_task_1

Висновок: При юзс **-erf** отримав більш валідні піки. Це через те, що в процесі навчання гранично випадкові ліси мають більше можливостей для вибору оптимальних дерев рішень, одже, як правило, вони забезпечують отримання кращих границь. Але кінцеві результати виявилися майже однаковими при використанні обох прапорців.

Завдання 2.2. Обробка дисбалансу класів.

```

1 import sys
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from sklearn.ensemble import ExtraTreesClassifier
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import classification_report
7 from utilities import visualize_classifier
8 # Завантаження вхідних даних
9 input_file = 'data_imbalance.txt'
10 data = np.loadtxt(input_file, delimiter=',')
11 X, y = data[:, :-1], data[:, -1]
12 # Поділ вхідних даних на два класи на підставі міток
13 class_0 = np.array(X[y == 0])
14 class_1 = np.array(X[y == 1])
15 # Візуалізація вхідних даних
16 plt.figure()
17 plt.scatter(class_0[:, 0], class_0[:, 1], s=75, facecolors='black',
18             edgecolors='black', linewidth=1, marker='x')
19 plt.scatter(class_1[:, 0], class_1[:, 1], s=75, facecolors='white',
20             edgecolors='black', linewidth=1, marker='o')
21 plt.title('Input data')
22 # Розбиття даних на навчальний та тестовий набори
23 X_train, X_test, y_train, y_test = train_test_split(
24     X, y, test_size=0.25, random_state=5)
25 # Класифікатор на основі гранично випадкових лісів
26 params = {'n_estimators': 100, 'max_depth': 4, 'random_state': 0}
27 if len(sys.argv) > 1:
28     if sys.argv[1] == 'balance':
29         params = {'n_estimators': 100, 'max_depth': 4, 'random_state': 0, 'class_weight': 'balanced'}
30     else:
31         raise TypeError("Invalid input argument: should be 'balance'")

```

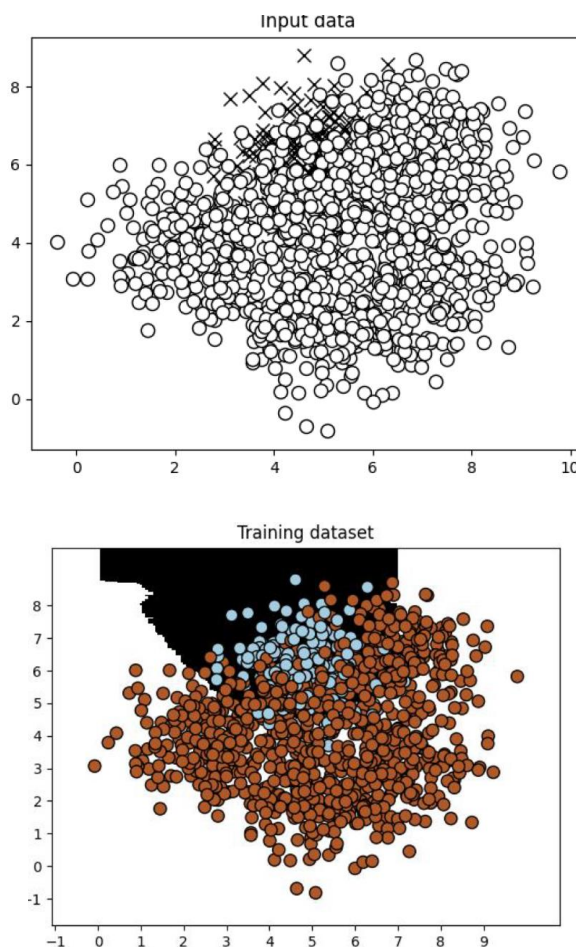
		Надворний М.Ю			ДУ «Житомирська політехніка».20.121.12 – Лр4	Арк.
		Філіпов В.О.				6
Змн.	Арк.	№ докум.	Підпис	Дата		


```

31         raise TypeError("Invalid input argument; should be 'balance'")
32     classifier = ExtraTreesClassifier(**params)
33     classifier.fit(X_train, y_train)
34     visualize_classifier(classifier, X_train, y_train, 'Training dataset')
35     y_test_pred = classifier.predict(X_test)
36     visualize_classifier(classifier, X_test, y_test, 'Test dataset')
37     # Обчислення показників ефективності класифікатора
38     class_names = ['Class-0', 'Class-1']
39     print("\n" + "#" * 40)
40     print("\nClassifier performance on training dataset\n")
41     print(classification_report(y_train, classifier.predict(X_train), target_names=class_names))
42     print("#" * 40 + "\n")
43     print("#" * 40)
44     print("\nClassifier performance on test dataset\n")
45     print(classification_report(y_test, y_test_pred, target_names=class_names))
46     print("#" * 40 + "\n")
47     plt.show()
48

```

Рис 3. Код програми файлу LR_4_task_2

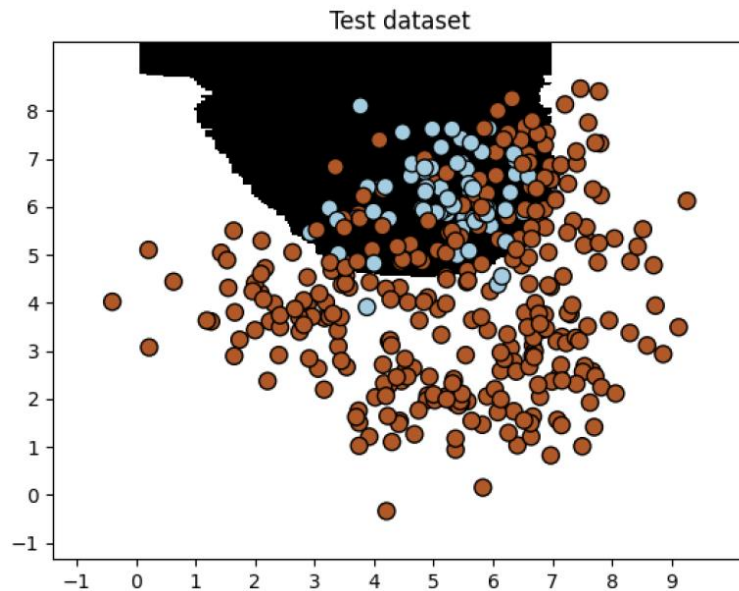


		Надворний М.Ю		
		Філіпов В.О.		
Змн.	Арк.	№ докум.	Підпис	Дата

ДУ «Житомирська політехніка».20.121.12 – Лр4

Арк.

7



Classifier performance on training dataset				
	precision	recall	f1-score	support
Class-0	0.44	0.93	0.60	181
Class-1	0.98	0.77	0.86	944
accuracy			0.78	375
macro avg	0.72	0.84	0.73	375
weighted avg	0.88	0.78	0.80	375

Рис 4. Результат виконання коду файлу LR_4_task_2

Завдання 2.3. Знаходження оптимальних навчальних параметрів за допомогою сіткового пошуку.


```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.metrics import classification_report
4 from sklearn.model_selection import GridSearchCV
5 from sklearn.ensemble import ExtraTreesClassifier
6 from sklearn.model_selection import train_test_split
7 from sklearn.metrics import classification_report
8 import pandas as pd
9 from utilities import visualize_classifier
10 input_file = 'data_random_forests.txt'
11 data = np.loadtxt(input_file, delimiter=',')
12 X, y = data[:, :-1], data[:, -1]
13 # Розбиття даних на три класи на підставі міток
14 class_0 = np.array(X[y == 0])
15 class_1 = np.array(X[y == 1])
16 class_2 = np.array(X[y == 2])
17 # Розбиття даних на навчальний та тестовий набори
18 X_train, X_test, y_train, y_test = train_test_split(
19     X, y, test_size=0.25, random_state=5)
20 # Визначення сітки значень параметрів
21 parameter_grid = [{'n_estimators': [100], 'max_depth': [2, 4, 7, 12, 16]},
22                   {'max_depth': [4], 'n_estimators': [25, 50, 100, 250]}]
23 ]
24 metrics = ['precision_weighted', 'recall_weighted']
25 for metric in metrics:
26     print("\n#### Searching optimal parameters for", metric)
27     classifier = GridSearchCV(
28         ExtraTreesClassifier(random_state=0),
29         parameter_grid, cv=5, scoring=metric)
30     classifier.fit(X_train, y_train)

```

```

classifier.fit(X_train, y_train)
df = pd.DataFrame(classifier.cv_results_)
df_columns_to_print = [column for column in df.columns if 'param' in column or 'score' in column]
print(df[df_columns_to_print])
print("\nBest parameters:", classifier.best_params_)
y_pred = classifier.predict(X_test)
print("\nPerformance report:\n")
print(classification_report(y_test, y_pred))

```

Рис 5. Код програми файлу LR_4_task_3

```

#### Searching optimal parameters for precision_weighted
mean_score_time  std_score_time  ... std_test_score rank_test_score
0      0.011244      3.890560e-04  ...      0.025132           1
1      0.011416      5.022505e-04  ...      0.023032           5
2      0.011253      4.321109e-04  ...      0.026560           4
3      0.011960      1.387637e-05  ...      0.028334           8
4      0.012765      7.462774e-04  ...      0.033429           9
5      0.003592      4.895782e-04  ...      0.029698           2
6      0.005984      9.536743e-08  ...      0.020401           7
7      0.010372      4.885192e-04  ...      0.023032           5
8      0.035558      1.915573e-02  ...      0.026867           3

[9 rows x 13 columns]

Best parameters: {'max_depth': 2, 'n_estimators': 100}

Performance report:

              precision    recall  f1-score   support

0.0         0.94         0.81         0.87         79
1.0         0.81         0.86         0.83         78
2.0         0.83         0.91         0.87         76

accuracy          0.86
macro avg          0.86         0.86         0.86         225
weighted avg          0.86         0.86         0.86         225

```

		Надворний М.Ю			ДУ «Житомирська політехніка».20.121.12 – Лр4	Арк.
		Філіпов В.О.				9
Змн.	Арк.	№ докум.	Підпис	Дата		

```
##### Searching optimal parameters for recall_weighted
| mean_score_time std_score_time ... std_test_score rank_test_score
0      0.012699      0.001604 ...      0.027075      1
1      0.014285      0.004264 ...      0.022468      5
2      0.014895      0.001292 ...      0.026749      3
3      0.013389      0.000439 ...      0.028497      8
4      0.017204      0.004969 ...      0.034744      9
5      0.004139      0.000478 ...      0.029407      1
6      0.007354      0.001175 ...      0.020096      7
7      0.012357      0.001032 ...      0.022468      5
8      0.034883      0.009096 ...      0.026749      3

[9 rows x 13 columns]

Best parameters: {'max_depth': 2, 'n_estimators': 100}

Performance report:

              precision    recall  f1-score   support

      0.0         0.94      0.81      0.87         79
      1.0         0.81      0.86      0.83         70
      2.0         0.83      0.91      0.87         76

 accuracy          0.86
 macro avg          0.86
weighted avg          0.86
```

Рис 6. Результат виконання коду файлу LR_4_task_3

Завдання 2.4. Обчислення відносної важливості ознак.

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.tree import DecisionTreeRegressor
4 from sklearn.ensemble import AdaBoostRegressor
5 from sklearn import datasets
6 from sklearn.metrics import mean_squared_error, explained_variance_score
7 from sklearn.model_selection import train_test_split
8 from sklearn.utils import shuffle
9 # Завантаження даних із цінами на нерухомість
10 housing_data = datasets.load_boston()
11 # Перемішування даних
12 X, y = shuffle(housing_data.data, housing_data.target, random_state=7)
13 # Розбиття даних на навчальний та тестовий набори
14 X_train, X_test, y_train, y_test = train_test_split(
15     X, y, test_size=0.2, random_state=7)
16 # Модель на основі пересіпа AdaBoost
17 regressor = AdaBoostRegressor(DecisionTreeRegressor(max_depth=4),
18                               n_estimators=400, random_state=7)
19 regressor.fit(X_train, y_train)
20 # Обчислення показників ефективності пересіпа AdaBoost
21 y_pred = regressor.predict(X_test)
22 mse = mean_squared_error(y_test, y_pred)
23 evs = explained_variance_score(y_test, y_pred)
24 print("\nADABOOST REGRESSOR")
25 print("Mean squared error =", round(mse, 2))
26 print("Explained variance score =", round(evs, 2))
27 # Вилучення важливості ознак
28 feature_importances = regressor.feature_importances_
29 feature_names = housing_data.feature_names
30 # Нормалізація значень важливості ознак
31 feature_importances = 100.0 * (feature_importances / max(feature_importances))
```

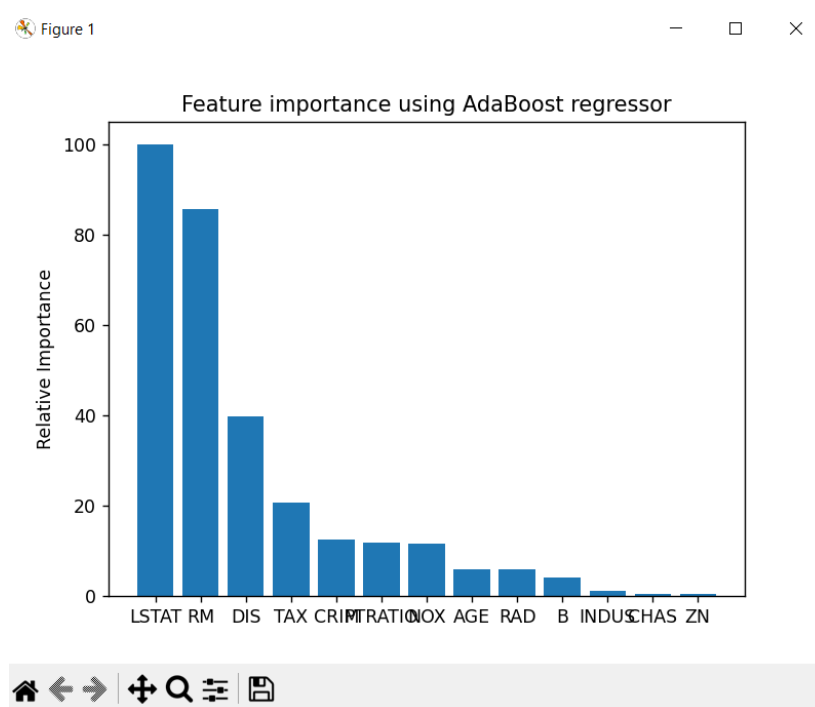
		Надворний М.Ю			ДУ «Житомирська політехніка».20.121.12 – Лр4	Арк.
		Філіпов В.О.				10
Змн.	Арк.	№ докум.	Підпис	Дата		

```

31 feature_importances = 100.0 * (feature_importances / max(feature_importances))
32 # Сортування та перестановка значень
33 index_sorted = np.flipud(np.argsort(feature_importances))
34 # Розміщення міток уздовж осі X
35 pos = np.arange(index_sorted.shape[0]) + 0.5
36 # Побудова стовпчастої діаграми
37 plt.figure()
38 plt.bar(pos, feature_importances[index_sorted], align='center')
39 plt.xticks(pos, feature_names[index_sorted])
40 plt.ylabel('Relative Importance')
41 plt.title('Feature importance using AdaBoost regressor')
42 plt.show()

```

Рис 7. Код програми файлу LR_4_task_4



```

ADABOOST REGRESSOR
Mean squared error = 22.7
Explained variance score = 0.79

Process finished with exit code 0

```

Рис 8. Результат виконання коду файлу LR_4_task_5

		Надворний М.Ю			ДУ «Житомирська політехніка».20.121.12 – Лр4	Арк.
		Філіпов В.О.				11
Змн.	Арк.	№ докум.	Підпис	Дата		

Завдання 2.5. Прогнозування інтенсивності дорожнього руху за допомогою класифікатора на основі гранично випадкових лісів.

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.metrics import classification_report, mean_absolute_error
4 from sklearn import preprocessing
5 from sklearn.model_selection import train_test_split
6 from sklearn.ensemble import ExtraTreesRegressor
7 from sklearn.metrics import classification_report
8 # Завантаження вхідних даних
9 input_file = 'traffic_data.txt'
10 data = []
11 with open(input_file, 'r') as f:
12     for line in f.readlines():
13         items = line[:-1].split(',')
14         data.append(items)
15 data = np.array(data)
16 # Перетворення рядкових даних на числові
17 label_encoder = []
18 X_encoded = np.empty(data.shape)
19 for i, item in enumerate(data[0]):
20     if item.isdigit():
21         X_encoded[:, i] = data[:, i]
22     else:
23         label_encoder.append(preprocessing.LabelEncoder())
24         X_encoded[:, i] = label_encoder[-1].fit_transform(data[:, i])
25 X = X_encoded[:, :-1].astype(int)
26 y = X_encoded[:, -1].astype(int)
27 # Розбиття даних на навчальний та тестовий набори
28 X_train, X_test, y_train, y_test = train_test_split(
29     X, y, test_size=0.25, random_state=5)
30 # Регресор на основі гранично випадкових лісів
31 params = {'n_estimators': 100, 'max_depth': 4, 'random_state': 0}
32 regressor = ExtraTreesRegressor(**params)
33 regressor.fit(X_train, y_train)
34 # Обчислення характеристик ефективності регресора на тестових даних
35 y_pred = regressor.predict(X_test)
36 print("Mean absolute error:", round(mean_absolute_error(y_test, y_pred), 2))
37 # Тестування кодування на одиночному прикладі
38 test_datapoint = ['Saturday', '10:20', 'Atlanta', 'no']
39 test_datapoint_encoded = [-1] * len(test_datapoint)
40 count = 0
41 for i, item in enumerate(test_datapoint):
42     if item.isdigit():
43         test_datapoint_encoded[i] = int(test_datapoint[i])
44     else:
45         test_datapoint_encoded[i] = int(label_encoder[count].transform([test_datapoint[i]]))
46         count = count + 1
47 test_datapoint_encoded = np.array(test_datapoint_encoded)
48 # Прогнозування результату для тестової точки даних
49 print("Predicted traffic:", int(regressor.predict([test_datapoint_encoded])[0]))

```

Рис 9. Код програми файлу LR_4_task_5

		Надворний М.Ю			ДУ «Житомирська політехніка».20.121.12 – Лр4	Арк.
		Філіпов В.О.				12
Змн.	Арк.	№ докум.	Підпис	Дата		

```
Mean absolute error: 7.42
Predicted traffic: 26

Process finished with exit code 0
```

Рис 10. Результат виконання коду файлу LR_4_task_5

Завдання 2.6. Створення навчального конвеєра (конвеєра машинного навчання).

```
1 from sklearn.datasets import _samples_generator
2 from sklearn.feature_selection import SelectKBest, f_regression
3 from sklearn.pipeline import Pipeline
4 from sklearn.ensemble import ExtraTreesClassifier
5 # Генерація даних
6 X, y = _samples_generator.make_classification(n_samples=150,
7                                             n_features=25, n_classes=3, n_informative=6,
8                                             n_redundant=0, random_state=7)
9 # Вибір k найважливіших ознак
10 k_best_selector = SelectKBest(f_regression, k=9)
11 # Ініціалізація класифікатора на основі гранично випадкового лісу
12 classifier = ExtraTreesClassifier(n_estimators=60, max_depth=4)
13 # Створення конвеєра
14 processor_pipeline = Pipeline([('selector', k_best_selector), ('erf', classifier)])
15 # Встановлення параметрів
16 processor_pipeline.set_params(selector__k=7, erf__n_estimators=30)
17 # Навчання конвеєра
18 processor_pipeline.fit(X, y)
19 # Прогнозування результатів для вхідних даних
20 output = processor_pipeline.predict(X)
21 print("\nPredicted output:\n", output)
22 # Виведення оцінки
23 print("\nScore:", processor_pipeline.score(X, y))
24 # Виведення ознак, відібраних селектором конвеєра
25 status = processor_pipeline.named_steps['selector'].get_support()
26 # Вилучення та виведення індексів обраних ознак
27 selected = [i for i, x in enumerate(status) if x]
28 print("\nIndices of selected features:", ', '.join([str(x) for x in selected]))
29
```

Рис 11. Код програми файлу LR_4_task_6

```
Predicted output:
[1 2 2 0 2 0 2 1 0 1 1 2 1 0 2 2 1 0 0 1 0 2 1 1 2 2 0 0 1 0 1 2 1 0 2 2 1
1 2 2 2 0 1 2 2 1 2 2 1 0 1 2 2 2 2 0 2 2 0 2 2 2 1 1 1 2 0 1 0 2
0 0 1 0 2 0 0 1 2 2 0 0 0 0 2 2 2 1 2 0 2 1 2 2 0 0 1 1 1 1 2 2 0 2 0 1 1
0 2 1 1 0 1 1 1 1 0 0 0 1 2 0 0 0 2 1 2 0 0 1 0 1 1 0 1 1 1 1 2 0 0 1 2 0
2 2]

Score: 0.8866666666666667

Indices of selected features: 4, 7, 8, 12, 14, 17, 22

Process finished with exit code 0
```

Рис 12. Результат виконання коду файлу LR_4_task_6

Висновок: Перший абзац містить прогнозовані вихідні мітки за допомогою конвеєра. Значення Score відображає ефективність конвеєра. Останній абзац містить індекси вибраних ознак.

Завдання 2.7. Пошук найближчих сусідів.

```

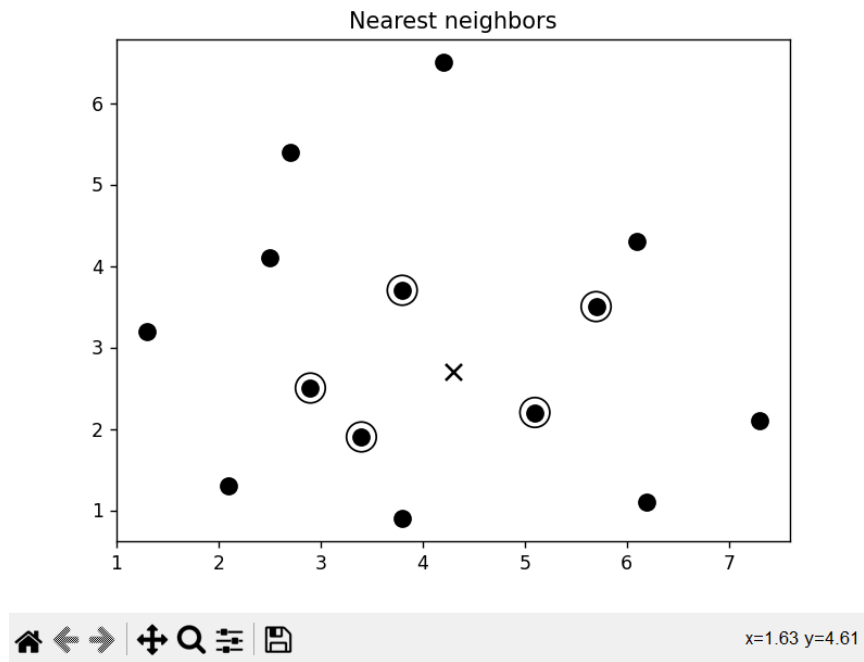
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.neighbors import NearestNeighbors
4 # Вхідні дані
5 X = np.array([[2.1, 1.3], [1.3, 3.2], [2.9, 2.5], [2.7, 5.4], [3.8, 0.9],
6              [7.3, 2.1], [4.2, 6.5], [3.8, 3.7], [2.5, 4.1], [3.4, 1.9],
7              [5.7, 3.5], [6.1, 4.3], [5.1, 2.2], [6.2, 1.1]])
8 # Кількість найближчих сусідів
9 k = 5
10 # Тестова точка даних
11 test_datapoint = [4.3, 2.7]
12 # Відображення вхідних даних на графіку
13 plt.figure()
14 plt.title('Input data')
15 plt.scatter(X[:, 0], X[:, 1], marker='o', s=75, color='black')
16 # Побудова моделі на основі методу k найближчих сусідів
17 knn_model = NearestNeighbors(n_neighbors=k, algorithm='ball_tree').fit(X)
18 distances, indices = knn_model.kneighbors([test_datapoint])
19 # Виведемо 'k' найближчих сусідів
20 print("\nK Nearest Neighbors:")
21 for rank, index in enumerate(indices[0][:k], start=1):
22     print(str(rank) + " ==> ", X[index])
23 # Візуалізація найближчих сусідів разом із тестовою точкою даних
24 plt.figure()
25 plt.title('Nearest neighbors')
26 plt.scatter(X[:, 0], X[:, 1], marker='o', s=75, color='k')
27 plt.scatter(X[indices[0][:k][:, 0], X[indices[0][:k][:, 1],
28             marker='o', s=250, color='k', facecolors='none')
29 plt.scatter(test_datapoint[0], test_datapoint[1],
30             marker='x', s=75, color='k')
31 plt.show()

```

Рис 13. Код програми файлу LR_4_task_7

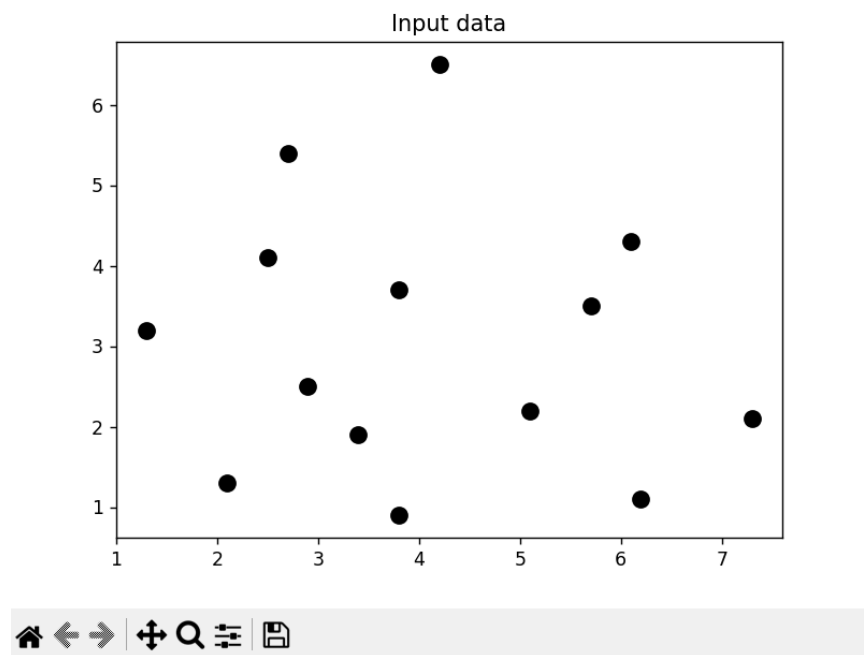
		Надворний М.Ю			ДУ «Житомирська політехніка».20.121.12 – Лр4	Арк.
		Філіпов В.О.				14
Змн.	Арк.	№ докум.	Підпис	Дата		

Figure 2



вхідні дані, тестову точку і її 5 найближчих сусідів

Figure 1



вхідні дані

		Надворний М.Ю.			ДУ «Житомирська політехніка».20.121.12 – Лр4	Арк.
		Філіпов В.О.				15
Змн.	Арк.	№ докум.	Підпис	Дата		

```

K Nearest Neighbors:
1 ==> [5.1 2.2]
2 ==> [3.8 3.7]
3 ==> [3.4 1.9]
4 ==> [2.9 2.5]
5 ==> [5.7 3.5]

Process finished with exit code 0

```

5 найближчих сусідів

Рис 14. Результат виконання коду файлу LR_4_task_7

Завдання 2.8. Створити класифікатор методом k найближчих сусідів.

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import matplotlib.cm as cm
4 from sklearn import neighbors, datasets
5 # Завантаження вхідних даних
6 input_file = 'data.txt'
7 data = np.loadtxt(input_file, delimiter=',')
8 X, y = data[:, :-1], data[:, -1].astype(np.int)
9 # Відображення вхідних даних на графіку
10 plt.figure()
11 plt.title('Input data')
12 marker_shapes = 'v^os'
13 mapper = [marker_shapes[i] for i in y]
14 for i in range(X.shape[0]):
15     plt.scatter(X[i, 0], X[i, 1], marker=mapper[i],
16               s=75, edgecolors='black', facecolors='none')
17 # Кількість найближчих сусідів
18 num_neighbors = 12
19 # Розмір кроку сітки візуалізації
20 step_size = 0.01
21 # Створення класифікатора на основі методу k найближчих сусідів
22 classifier = neighbors.KNeighborsClassifier(num_neighbors, weights='distance')
23 # Навчання моделі на основі методу k найближчих сусідів
24 classifier.fit(X, y)
25 # Створення сітки для відображення меж на графіку
26 x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
27 y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
28 x_values, y_values = np.meshgrid(np.arange(x_min, x_max, step_size), np.arange(y_min, y_max, step_size))
29 # Виконання класифікатора на всіх точках сітки
30 output = classifier.predict(np.c_[x_values.ravel(), y_values.ravel()])
31 # Візуалізація результатів

```

		Надворний М.Ю			ДУ «Житомирська політехніка».20.121.12 – Лр4	Арк.
		Філіпов В.О.				16
Змн.	Арк.	№ докум.	Підпис	Дата		

```

31 # Візуалізація передбачуваного результату
32 output = output.reshape(x_values.shape)
33 plt.figure()
34 plt.pcolormesh(x_values, y_values, output, cmap=cm.Paired)
35 # Накладання навчальних точок на карту
36 for i in range(X.shape[0]):
37     plt.scatter(X[i, 0], X[i, 1], marker=mapper[i],
38                 s=50, edgecolors='black', facecolors='none')
39 plt.xlim(x_values.min(), x_values.max())
40 plt.ylim(y_values.min(), y_values.max())
41 plt.title('K Nearest Neighbors classifier model boundaries')
42 # Тестування вхідної точки даних
43 test_datapoint = [5.1, 3.6]
44 plt.figure()
45 plt.title('Test datapoint')
46 for i in range(X.shape[0]):
47     plt.scatter(X[i, 0], X[i, 1], marker=mapper[i],
48                 s=75, edgecolors='black', facecolors='none')
49 plt.scatter(test_datapoint[0], test_datapoint[1], marker='x',
50             linewidth=6, s=200, facecolors='black')
51 # Вилучення K найближчих сусідів
52 _, indices = classifier.kneighbors([test_datapoint])
53 indices = indices.astype(np.int)[0]
54 # Відображення K найближчих сусідів на графіку
55 plt.figure()
56 plt.title('K Nearest Neighbors')
57 for i in indices:
58     plt.scatter(X[i, 0], X[i, 1], marker=mapper[y[i]],
59                 linewidth=3, s=100, facecolors='black')
60 plt.scatter(test_datapoint[0], test_datapoint[1], marker='x',
61             linewidth=6, s=200, facecolors='black')

```

```

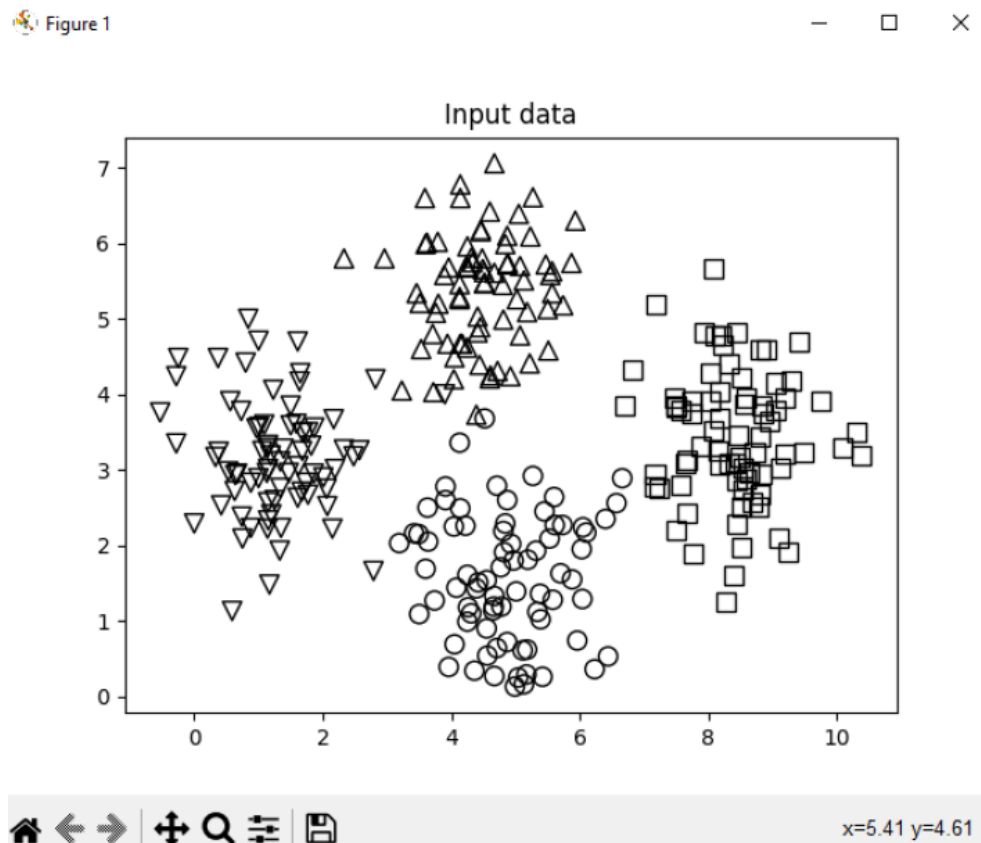
54 # Відображення K найближчих сусідів на графіку
55 plt.figure()
56 plt.title('K Nearest Neighbors')
57 for i in indices:
58     plt.scatter(X[i, 0], X[i, 1], marker=mapper[y[i]],
59                 linewidth=3, s=100, facecolors='black')
60 plt.scatter(test_datapoint[0], test_datapoint[1], marker='x',
61             linewidth=6, s=200, facecolors='black')
62 for i in range(X.shape[0]):
63     plt.scatter(X[i, 0], X[i, 1], marker=mapper[i],
64                 s=75, edgecolors='black', facecolors='none')
65 print("Predicted output:", classifier.predict([test_datapoint])[0])
66 plt.show()
67

```

Рис 15. Код програми файлу LR_4_task_8

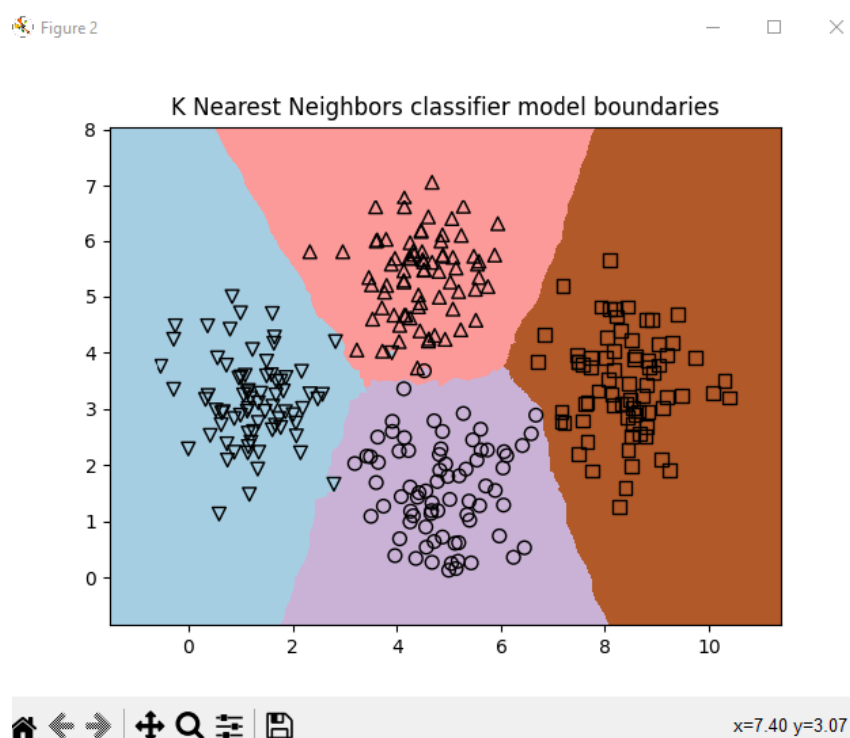
		Надворний М.Ю			ДУ «Житомирська політехніка».20.121.12 – Лр4	Арк.
		Філіпов В.О.				17
Змн.	Арк.	№ докум.	Підпис	Дата		

Figure 1



вхідні дані

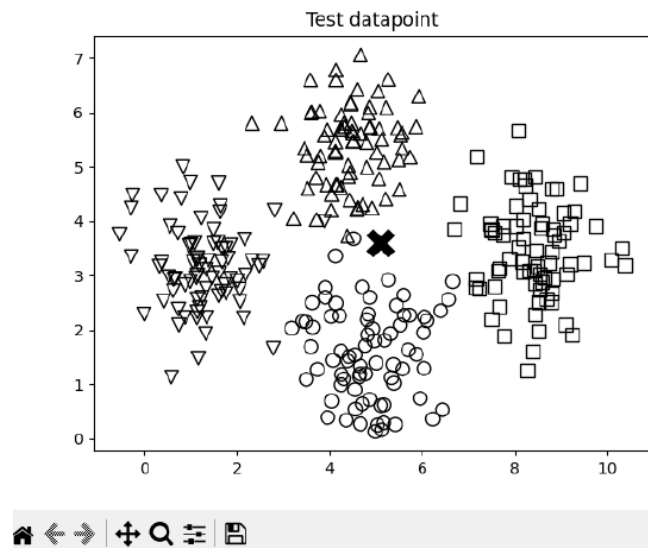
Figure 2



межі класифікатора

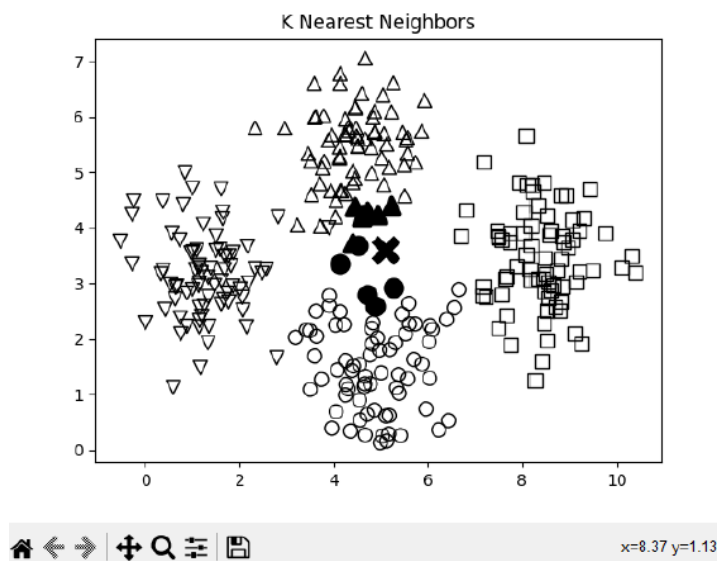
		Надворний М.Ю			ДУ «Житомирська політехніка».20.121.12 – Лр4	Арк.
		Філіпов В.О.				18
Змн.	Арк.	№ докум.	Підпис	Дата		

Figure 3



тестова точка до вхідного набору даних

Figure 4



12 найближчих сусідів

Рис 16. Результат виконання коду файлу LR_4_task_8

Завдання 2.9. Обчислення оцінок подібності.

		Надворний М.Ю			ДУ «Житомирська політехніка».20.121.12 – Лр4	Арк.
		Філіпов В.О.				19
Змн.	Арк.	№ докум.	Підпис	Дата		

```

1 import argparse
2 import json
3 import numpy as np
4 def build_arg_parser():
5     parser = argparse.ArgumentParser(description='Compute similarity score')
6     parser.add_argument('--user1', dest='user1', required=True, help='First user')
7     parser.add_argument('--user2', dest='user2', required=True,
8                           help='Second user')
9     parser.add_argument('--score-type", dest="score_type", required=True,
10                        choices=['Euclidean', 'Pearson'], help='Similarity metric to be used')
11     return parser
12 # Обчислення оцінки евклідова відстані між користувачами user1 та user2
13 def euclidean_score(dataset, user1, user2):
14     if user1 not in dataset:
15         raise TypeError('Cannot find ' + user1 + ' in the dataset')
16
17     if user2 not in dataset:
18         raise TypeError('Cannot find ' + user2 + ' in the dataset')
19     # Фільми, оцінені обома користувачами, user1 та user2
20     common_movies = {}
21     for item in dataset[user1]:
22         if item in dataset[user2]:
23             common_movies[item] = 1
24     # За відсутності фільмів, оцінених обома користувачами, оцінка приймається рівною 0
25     if len(common_movies) == 0:
26         return 0
27     squared_diff = []
28     for item in dataset[user1]:
29         if item in dataset[user2]:
30             squared_diff.append(np.square(dataset[user1][item] - dataset[user2][item]))
31     return 1 / (1 + np.sqrt(np.sum(squared_diff)))

```

```

32 # Обчислення кореляційної оцінки Пірсона між користувачем1 і користувачем2
33 def pearson_score(dataset, user1, user2):
34     if user1 not in dataset:
35         raise TypeError('Cannot find ' + user1 + ' in the dataset')
36     if user2 not in dataset:
37         raise TypeError('Cannot find ' + user2 + ' in the dataset')
38     # Фільми, оцінені обома користувачами, user1 та user2
39     common_movies = {}
40     for item in dataset[user1]:
41         if item in dataset[user2]:
42             common_movies[item] = 1
43     num_ratings = len(common_movies)
44     # За відсутності фільмів, оцінених обома користувачами, оцінка приймається рівною 0
45     if num_ratings == 0:
46         return 0
47     # Обчислення суми рейтингових оцінок усіх фільмів, оцінених обома користувачами
48     user1_sum = np.sum([dataset[user1][item] for item in common_movies])
49     user2_sum = np.sum([dataset[user2][item] for item in common_movies])
50     # Обчислення суми квадратів рейтингових оцінок всіх фільмів, оцінених обома користувачами
51     user1_squared_sum = np.sum([np.square(dataset[user1][item]) for item in common_movies])
52     user2_squared_sum = np.sum([np.square(dataset[user2][item]) for item in common_movies])
53     # Обчислення суми творів рейтингових оцінок всіх фільмів, оцінених обома користувачами
54     sum_of_products = np.sum([dataset[user1][item] * dataset[user2][item] for item in common_movies])
55     # Обчислення коефіцієнта кореляції Пірсона
56     Sxy = sum_of_products - (user1_sum * user2_sum / num_ratings)
57     Sxx = user1_squared_sum - np.square(user1_sum) / num_ratings
58     Syy = user2_squared_sum - np.square(user2_sum) / num_ratings
59     if Sxx * Syy == 0:
60         return 0
61     return Sxy / np.sqrt(Sxx * Syy)

```

		Надворний М.Ю			ДУ «Житомирська політехніка».20.121.12 – Лр4	Арк.
		Філіпов В.О.				20
Змн.	Арк.	№ докум.	Підпис	Дата		


```

62 if __name__ == '__main__':
63     args = build_arg_parser().parse_args()
64     user1 = args.user1
65     user2 = args.user2
66     score_type = args.score_type
67     ratings_file = 'ratings.json'
68     with open(ratings_file, 'r') as f:
69         data = json.loads(f.read())
70     if score_type == 'Euclidean':
71         print("\nEuclidean score:")
72         print(euclidean_score(data, user1, user2))
73     else:
74         print("\nPearson score:")
75         print(pearson_score(data, user1, user2))
76

```

Рис 17. Код програми файлу LR_4_task_9

```

Pearson score:
.9909924304103233
$ C:\Users\Admin\PycharmProjects\labka4> python LR_4_task_9.py --user1 "David Smith" --user2 "Bill Duffy" --score-type Euclidean

Euclidean score:
.585786437626905
$ C:\Users\Admin\PycharmProjects\labka4> python LR_4_task_9.py --user1 "David Smith" --user2 "Bill Duffy" --score-type Pearson

Pearson score:
.9909924304103233
$ C:\Users\Admin\PycharmProjects\labka4>

```

```

Pearson score:
-0.723679961015113
PS C:\Users\Admin\PycharmProjects\labka4> python LR_4_task_9.py --user1 "David Smith" --user2 "Samuel Miller" --score-type Euclidean

Euclidean score:
0.30383243470068705
PS C:\Users\Admin\PycharmProjects\labka4> python LR_4_task_9.py --user1 "David Smith" --user2 "Samuel Miller" --score-type Pearson

Pearson score:
0.7587869100593281
PS C:\Users\Admin\PycharmProjects\labka4> python LR_4_task_9.py --user1 "David Smith" --user2 "Julie Hanneel" --score-type Euclidean

Euclidean score:
0.2857142857142857
PS C:\Users\Admin\PycharmProjects\labka4> python LR_4_task_9.py --user1 "David Smith" --user2 "Julie Hanneel" --score-type Pearson

Pearson score:
0
PS C:\Users\Admin\PycharmProjects\labka4> python LR_4_task_9.py --user1 "David Smith" --user2 "Clarissa Jackson" --score-type Euclidean

Euclidean score:
Traceback (most recent call last):
  File "C:\Users\Admin\PycharmProjects\labka4\LR_4_task_9.py", line 100, in <module>
    print(euclidean_score(data, user1, user2))
  File "C:\Users\Admin\PycharmProjects\labka4\LR_4_task_9.py", line 22, in euclidean_score
    raise ValueError('Cannot find " + user2 + " in the dataset')
ValueError: Cannot find Clarissa Jackson in the dataset
PS C:\Users\Admin\PycharmProjects\labka4> python LR_4_task_9.py --user1 "David Smith" --user2 "Clarissa Jackson" --score-type Euclidean

Euclidean score:
0.2898979485563564
PS C:\Users\Admin\PycharmProjects\labka4> python LR_4_task_9.py --user1 "David Smith" --user2 "Clarissa Jackson" --score-type Pearson
usage: LR_4_task_9.py [-h] --user1 USER1 --user2 USER2 --score-type {Euclidean,Pearson}
LR_4_task_9.py: error: argument --score-type: invalid choice: 'Pearson' (choose from 'Euclidean', 'Pearson')
PS C:\Users\Admin\PycharmProjects\labka4> python LR_4_task_9.py --user1 "David Smith" --user2 "Clarissa Jackson" --score-type Pearson

Pearson score:
0.6944217042199275
PS C:\Users\Admin\PycharmProjects\labka4> python LR_4_task_9.py --user1 "David Smith" --user2 "Adam Cohen" --score-type Euclidean

Euclidean score:
0.38742588672279384
PS C:\Users\Admin\PycharmProjects\labka4> python LR_4_task_9.py --user1 "David Smith" --user2 "Adam Cohen" --score-type Pearson

Pearson score:
0.9881882718950217
PS C:\Users\Admin\PycharmProjects\labka4> python LR_4_task_9.py --user1 "David Smith" --user2 "Chris Duncan" --score-type Euclidean

Euclidean score:
0.38742588672279384
PS C:\Users\Admin\PycharmProjects\labka4> python LR_4_task_9.py --user1 "David Smith" --user2 "Chris Duncan" --score-type Pearson

Pearson score:
1.0
PS C:\Users\Admin\PycharmProjects\labka4>

```

Рис 18. Результат виконання коду файлу LR_4_task_9

		Надворний М.Ю			ДУ «Житомирська політехніка».20.12.12 – Лр4	Арк.
		Філіпов В.О.				21
Змн.	Арк.	№ докум.	Підпис	Дата		

Оцінка подібності за Пірсоном демонструє кращі результати в порівнянні з евклідовою оцінкою подібності.

Завдання 2.10. Пошук користувачів зі схожими уподобаннями методом колаборативної фільтрації.

```

1 import argparse
2 import json
3 import numpy as np
4 from LR_4_task_9m import pearson_score
5 def build_arg_parser():
6     parser = argparse.ArgumentParser(description='Find users who are similar to the in-put user')
7     parser.add_argument('--user', dest='user', required=True,
8                           help='Input user')
9     return parser
10 # Знаходження користувачів у наборі даних, схожих на введеного користувача
11 def find_similar_users(dataset, user, num_users):
12     if user not in dataset:
13         raise TypeError('Cannot find ' + user + ' in the dataset')
14     # Обчислення оцінки подібності за Пірсоном між
15     # вказаним користувачем та всіма іншими
16     # користувачами в наборі даних
17     scores = np.array([[x, pearson_score(dataset, user,
18                                         x)] for x in dataset if x != user])
19     # Сортування оцінок за спаданням
20     scores_sorted = np.argsort(scores[:, 1])[::-1]
21     # Вилучення оцінок перших 'num_users' користувачів
22     top_users = scores_sorted[:num_users]
23     return scores[top_users]
24 if __name__ == '__main__':
25     args = build_arg_parser().parse_args()
26     user = args.user
27     ratings_file = 'ratings.json'
28     with open(ratings_file, 'r') as f:
29         data = json.loads(f.read())
30     print('\nUsers similar to ' + user + ':\n')
31     similar_users = find_similar_users(data, user, 3)
32     print('\nUsers similar to ' + user + ':\n')
33     similar_users = find_similar_users(data, user, 3)
34     print('User\t\t\tSimilarity score')
35     print('-' * 41)
36     for item in similar_users:
37         print(item[0], '\t\t', round(float(item[1]), 2))

```

Рис 19. Код програми файлу LR_4_task_10

		Надворний М.Ю			ДУ «Житомирська політехніка».20.121.12 – Лр4	Арк.
		Філіпов В.О.				22
Змн.	Арк.	№ докум.	Підпис	Дата		

```

LR_4_task_10.py: error: the following arguments are required: --user
PS C:\Users\Admin\PycharmProjects\labka4> python LR_4_task_10.py --user "Clarissa Jackson"

Users similar to Clarissa Jackson:

User                Similarity score
-----
Chris Duncan        1.0
Bill Duffy           0.83
Samuel Miller        0.73
PS C:\Users\Admin\PycharmProjects\labka4> python LR_4_task_10.py --user "Bill Duffy"

Users similar to Bill Duffy:

User                Similarity score
-----
David Smith          0.99
Samuel Miller         0.88
Adam Cohen            0.86
PS C:\Users\Admin\PycharmProjects\labka4>

```

Рис 20. Результат виконання коду файлу LR_4_task_10

Юзер “Clarissa Jackson” має однакові вподобання з користувачем “Chris Duncan”, а користувач “Bill Duffy” – майже однакові з “David Smith”.

Завдання 2.11. Створення рекомендаційної системи фільмів.

```

1  import argparse
2  import json
3  import numpy as np
4  from LR_4_task_9m import pearson_score
5  from LR_4_task_10m import find_similar_users
6  def build_arg_parser():
7      parser = argparse.ArgumentParser(description='Find the movie recommendations for the given user')
8      parser.add_argument('--user', dest='user', required=True,
9                          help='Input user')
10     return parser
11     # Отримання рекомендації щодо фільмів для вказаного користувача
12     def get_recommendations(dataset, input_user):
13         if input_user not in dataset:
14             raise TypeError('Cannot find ' + input_user + ' in the dataset')
15         overall_scores = {}
16         similarity_scores = {}
17         for user in [x for x in dataset if x != input_user]:
18             similarity_score = pearson_score(dataset, input_user, user)
19             if similarity_score <= 0:
20                 continue
21             filtered_list = [x for x in dataset[user] if x not in \
22                             dataset[input_user] or dataset[input_user][x] == 0]
23             for item in filtered_list:
24                 overall_scores.update({item: dataset[user][item] * similarity_score})
25                 similarity_scores.update({item: similarity_score})
26         if len(overall_scores) == 0:
27             return ['No recommendations possible']
28         # Генерація рейтингів фільмів за допомогою їх нормалізації
29         movie_scores = np.array([[score / similarity_scores[item], item]
30                                   for item, score in overall_scores.items()])

```

		Надворний М.Ю			ДУ «Житомирська політехніка».20.121.12 – Лр4	Арк.
		Філіпов В.О.				23
Змн.	Арк.	№ докум.	Підпис	Дата		

```

34     movie_recommendations = [movie for _, movie in movie_scores]
35     return movie_recommendations
36 ► if __name__ == '__main__':
37     args = build_arg_parser().parse_args()
38     user = args.user
39     ratings_file = 'ratings.json'
40     with open(ratings_file, 'r') as f:
41         data = json.loads(f.read())
42     print("\nMovie recommendations for " + user + ":")
43     movies = get_recommendations(data, user)
44     for i, movie in enumerate(movies):
45         print(str(i + 1) + '. ' + movie)

```

Рис 21. Код програми файлу LR_4_task_11

```

PS C:\Users\Admin\PycharmProjects\labka4> python LR_4_task_11.py --user "Julie Hammel"

Movie recommendations for Julie Hammel:
1. The Apartment
2. Vertigo
3. Raging Bull
PS C:\Users\Admin\PycharmProjects\labka4> python LR_4_task_11.py --user "Clarissa Jackson"

Movie recommendations for Clarissa Jackson:
1. No recommendations possible
PS C:\Users\Admin\PycharmProjects\labka4> 

```

Рис 22. Результат виконання коду файлу LR_4_task_11

Julie Hammel = 3 реки

Кларіси = 0

Висновок: в даній лабораторній роботі навчився використовувати спеціалізовані бібліотеки та мову програмування Python, досліджувати методи ансамблів у машинному навчанні та створювати рекомендаційні системи.

Посилання на гіт: <https://github.com/Max1648/Artificial-Intelligence2>.

		Надворний М.Ю			ДУ «Житомирська політехніка».20.121.12 – Лр4	Арк.
		Філіпов В.О.				24
Змн.	Арк.	№ докум.	Підпис	Дата		